MICROCOPY RESOLUTION TEST CHART

NATIONAL BUREAU OF STANDARDS-1963-A

AFOSR-TR- ∴ ∴ 0 8 0 4

<span style="writing-mode: vertical">AD-A160 301</span>

DIFFERENTIAL METRICS IN PROBABILITY
SPACES BASED ON ENTROPY AND DIVERGENCE
MEASURES*


By


C. Radhakrishna Rao
University of Pittsburgh
U.S.A.


April 1985


Technical Report No. 85-08


<span style="writing-mode: vertical">DTIC FILE COPY</span>

Center for Multivariate Analysis
Fifth Floor Thackeray Hall
University of Pittsburgh
Pittsburgh, PA 15260

DTIC
ELECTE
OCT 1 5 1985
E

85 10 11 120

# DIFFERENTIAL METRICS IN PROBABILITY
## SPACES BASED ON ENTROPY AND DIVERGENCE
### MEASURES

C. Radhakrishna Rao

SUMMARY: In this paper are discussed some general methods of metrizing prob-
ability spaces through the introduction of a quadratic differential metric
in the parameter manifold of a set of probability distributions. These methods extend
the investigation made in Rao (1945) where the Fisher information matrix was
used to construct the metric, and the geodesic distance was suggested as a mea-
sure of dissimilarity between probability distributions.

The basic approach in the present *this* paper is first to construct a divergence
or a dissimilarity measure between any two probability distributions, and use it
to derive a differential metric by considering two distributions whose characterizing
parameters are close to each other. One measure of divergence considered is the
Jensen difference based on an entropy functional as defined in Rao (1982a).
Another is the f-divergence measure studied by Csiszár (1967). The latter
class leads to the differential metric based on the Fisher information matrix.
The geodesic distances based on this metric computed by various authors are
listed.

KEY WORDS: cross entropy, f-divergence, geodesic distance, information matrix,
Jensen difference, quadratic entropy.

# 1. INTRODUCTION

In an early paper (Rao, 1945), the author introduced a Riemannian (quadratic differential) metric over the space of a parametric family of probability distributions and proposed the geodesic distance induced by the metric as a measure of dissimilarity between probability distributions. The metric was based on the Fisher information matrix and it arose in a natural way through the concepts of statistical discrimination (Rao, 1949, 1954, 1973 pp. 329-332, 1982a). Such a choice of the quadratic differential metric, which we will refer to as the information metric, has indeed some attractive properties such as invariance for transformation of the variables as well as the parameters. It also seems to provide an appropriate (informative) geometry on the probability space for studying large sample properties of estimators of parameters in terms of simple loss functions as demonstrated by Amari (1982, 1983), Čencov (1982), Efron (1975, 1982), Eguchi (1983, 1984) and others.

The geodesic distances based on the information metric have been computed for a number of parametric family of distributions in recent papers by Atkinson and Mitchell (1981), Burbea (1984), Mitchell and Krzanowski (1985), and Oller and Cuadras (1985).

In two papers, Burbea and Rao (1982a, 1982b) gave some general methods for constructing quadratic differential metrics on probability spaces, of which the Fisher information metric belonged to a special class. In view of the rich variety of possible metrics, it would be useful to lay down some criteria for the choice of an appropriate metric for a given problem. Amari has stated that a metric should reflect the stochastic and statistical properties of the family of probability distributions. In particular he emphasized the invariance

of the metric under transformations of the variables as well as the parameters.
Čencov (1972) shows that the Fisher information metric is unique under some conditions including invariance. Burbea and Rao (1982a) showed that the Fisher information metric is the only metric associated with invariant divergence measures of the type introduced by Ciszàr (1967). However, there exist other types of invariant metrics as shown in Section 3 of this paper.

The choice of a metric naturally depends on a particular problem under investigation, and invariance may or may not be relevant. For instance, consider the space of multinomial distributions, $\Delta = \{(p_1,\ldots,p_n): p_i > 0, \Sigma p_i = 1\}$, which is a submanifold of the positive orthant, $X = \{(x_1,\ldots,x_n): x_i > 0\}$ of the Euclidean space $R^n$. A Riemannian metric on X automatically provides a metric on the submanifold $\Delta$. In a study of linkage and selection of gametes in a biological population, Shahshahani (1979) considered the metric

$$(1.1) \qquad ds^2 = \sum_1^n \frac{\Sigma x_i}{x_i} \, dx_i^2$$

which induces the information metric on $\Delta$. This metric provided a convenient framework for a discussion of certain biological problems. However, Nei (1978) considered a distance measure associated with the Euclidean metric

$$(1.2) \qquad ds^2 = \Sigma dx_i^2$$

which he found to be more appropriate for evolutionary studies in biology. The metric induced on $\Delta$ by (1.2) is not the Fisher information metric. Rao (1982a, 1982b) has shown that a more general type of metric

$$(1.3) \qquad \Sigma\Sigma a_{ij} dp_i dp_j$$

on $\Delta$, called the quadratic entropy is more meaningful in certain sociometric and

biometric studies.

The object of the present paper is to provide some general methods of constructing Riemannian metrics on probability spaces, and discuss in particular the metric generated by the quadratic entropy which is an ideal measure of diversity (see Lau, 1985 and Rao, 1982b), and has properties similar to the information metric, like invariance. We also give a list of geodesic distances based on the information metric computed by various authors (Atkinson and Mitchell, 1981; Burbea, 1984; Mitchell and Krzanowski, 1985; Oller and Cuadras, 1985 and Rao, 1945).

The basic approach adopted in the paper is first to define a measure of divergence or dissimilarity between two probability measures, and use it to derive a metric on M, the manifold of parameters, by considering two distributions defined by two contiguous points in M. We thus provide a wider basis for the construction of an appropriate geometry or geometries on the parameter space for discussion of practical problems. Some divergence measures may be more appropriate for discussing properties of estimators using simple loss functions while others may be appropriate in the study of population dynamics in biology. It is not unusual in practice to study a problem under different models for observed data to examine consistency and robustness of results. The variety of metrics reported in the paper would be of some use in this direction.

## 2. JENSEN DIFFERENCE AND ENTROPY DIFFERENTIAL METRIC

Let $\nu$ be a $\sigma$-finite additive measure defined on a $\sigma$-algebra of subsets of a measurable space $X$, and $P$ be the usual Lebesgue space of $\nu$ measurable density functions,

$$(2.1) \qquad P = \{p(x): \quad p(x) \geq 0, \ x \in X, \ \int_X p(x)d\nu(x) = 1\}.$$

We call H: $P \to R$ an entropy (functional) on $P$ if

(i) $H(p) = 0$ when p is degenerate,

(ii) $H(p)$ is concave on $P$.

In such a case, with $\lambda \geq 0$, $\mu \geq 0$, $\lambda + \mu = 1$, Rao (1982a) defined the Jensen difference between p and $q \in P$ as

$$(2.2) \qquad J(\lambda, \mu; \ p, q) = H(\lambda p + \mu q) - \lambda H(p) - \mu H(q).$$

The function J: $P \times P \to R$ is non-negative and vanishes if $p = q$ (iff $p = q$ when H is strictly concave). If the entropy function H is regarded as a measure of diversity within a population, then the Jensen difference J can be interpreted as a measure of diversity (or dissimilarity) between two populations. For the use of Jensen difference in the measurement, apportionment and analysis of diversity between populations, the reader is referred to Rao (1982a, 1982b).

Let us now consider a subset of probability densities characterized by a vector parameter $\theta$

$$P_\theta = \{p(x, \theta): \quad p(x, \theta) \in P, \ \theta \in M, \ \text{a manifold in } R^n\}$$

and assume that $p(x, \theta)$ is a smooth function admitting derivatives of a certain order with respect to $\theta$ and differention under the integral sign. For convenience of notation, we write

$$p(\cdot, \theta) = p_\theta, \quad H(\theta) = H(p_\theta), \quad H(\theta, \phi) = H(\lambda p_\theta + \mu p_\phi)$$

$$(2.3) \qquad J(\theta, \phi) = H(\theta, \phi) - \lambda H(\theta) - \mu H(\phi)$$

where $\theta, \phi \in M$. Putting $\phi = \theta + d\theta$ and denoting the i-th component of a vector with a subscript i, we consider the formal expansion of $J(\theta, \theta + d\theta)$,

$$(2.4) \quad \frac{1}{2!} \sum_{1}^{n}\sum_{1}^{n} \frac{\partial^2 J(\theta, \phi = \theta)}{\partial \phi_i \partial \phi_j} d\theta_i d\theta_j + \frac{1}{3!} \sum_{1}^{n}\sum_{1}^{n}\sum_{1}^{n} \frac{\partial^3 J(\theta, \phi = \theta)}{\partial \phi_i \partial \phi_j \partial \phi_k} d\theta_i d\theta_j d\theta_k + \ldots$$

$$= \frac{1}{2!} \Sigma\Sigma \, g_{ij}^H(\theta) d\theta_i d\theta_j + \frac{1}{3!} \Sigma\Sigma\Sigma \, c_{ijk}^H(\theta) d\theta_i d\theta_j d\theta_k + \ldots$$

In (2.4), the coefficients of the first order differentials vanish since $J(\theta, \phi)$ has a minimum at $\phi = \theta$, and the notation such as $\partial^2 J(\theta, \phi = \theta)/\partial \phi_i \partial \phi_j$ is used for replacing $\phi$ by $\theta$ after carrying out the indicated differentiations.

From the definition of the J function, it follows that the $(g_{ij}^H)$ is a non-negative definite matrix and obeys the tensorial law under transformation of parameters. We define the matrix and the associated metric

$$(2.5) \quad (g_{ij}^H) \text{ and } \Sigma\Sigma \, g_{ij}^H d\theta_i d\theta_j$$

as the H-entropy information matrix and H-entropy differential metric respectively. We prove the following theorem which provides an alternative computation of the H-information matrix directly from a given entropy H.

Theorem 2.1

$$(2.6) \quad g_{ij}^H(\theta) = - \frac{\partial^2 H(\lambda p_\theta + \mu p_\phi)}{\partial \theta_i \partial \phi_j} \bigg|_{\phi = \theta}.$$

Proof: By definition

$$g_{ij}^H(\theta) = \frac{\partial^2 J(\theta, \phi = \theta)}{\partial \phi_i \partial \phi_j}$$

$$(2.7) \qquad = \frac{\partial^2 H(\theta, \phi = \theta)}{\partial \phi_i \partial \phi_j} - \mu \frac{\partial^2 H(\phi = \theta)}{\partial \phi_i \partial \phi_j}.$$

Since $J(\theta, \phi)$ attains a minimum at $\phi = \theta$

$$(2.8) \qquad \frac{\partial H(\theta,\phi=\theta)}{\partial \phi_j} = \mu \frac{\partial H(\theta)}{\partial \theta_j} .$$

Differentiating both sides of (2.8) with respect to $\theta_i$ we have

$$(2.9) \qquad \frac{\partial^2 H(\theta,\phi=\theta)}{\partial \theta_i \partial \phi_j} + \frac{\partial^2 H(\theta,\phi=\theta)}{\partial \phi_i \partial \phi_j} = \mu \frac{\partial^2 H(\theta)}{\partial \theta_i \partial \theta_j}$$

which gives (2.6), and the desired result is proved.

Let us consider a general entropy function of the type

$$(2.10) \qquad H(p_\theta) = - \int h(p_\theta) d\nu(x)$$

where $h''$ is a non-negative function. Then using (2.6)

$$(2.11) \qquad g_{ij}^H(\theta) = g_{ij}^h(\theta) = - \frac{\partial^2 H(\theta,\phi=\theta)}{\partial \theta_i \partial \phi_j}$$

$$= \int \frac{\partial^2 h(\lambda p_\theta + \mu p_\phi)}{\partial \theta_i \partial \phi_j} \Big|_{\phi=\theta} d\nu(x)$$

$$= \lambda\mu \int h''(p_\theta) \frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x) .$$

If $[-h(x)] = -x \log x$, leading to Shannon's entropy, then

$$(2.12) \qquad g_{ij}^h = g_{ij}(\theta) = \lambda\mu \int \frac{1}{p_\theta} \frac{\partial p_\theta}{\partial \theta_j} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x)$$

become the elements of Fisher's information matrix. If $h(x) = (\alpha-1)^{-1}(x^\alpha - x)$, $\alpha \neq 1$, we have the $\alpha$-order entropy of Havrda and Charvát and

$$(2.13) \qquad g_{ij}^h = g_{ij}^\alpha(\theta) = \alpha\lambda\mu \int p^\alpha \frac{\partial \log p}{\partial \theta_i} \frac{\partial \log p}{\partial \theta_j} d\nu(x)$$

which provide the elements of $\alpha$-order entropy information matrix, and the corresponding differential metric given in Burbea and Rao (1982a, 1982b).

We prove Theorem 2.2 which gives alternative expressions for the coefficients of the third order differentials in the expansion of $J(\theta,\phi)$.

<u>Theorem 2.2</u>

$$(2.14) \quad c_{ijk}^{H} = -[\frac{\partial^3 H(\theta,\phi=\theta)}{\partial\theta_i\partial\theta_j\partial\phi_k} + \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\theta_i\partial\phi_j\partial\phi_k} + \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\theta_j\partial\phi_i\partial\phi_k}].$$

<u>Proof:</u>  By definition

$$(2.15) \quad c_{ijk}^{H}(\theta) = \frac{\partial^3 J(\theta,\phi=\theta)}{\partial\phi_i\partial\phi_j\partial\phi_k}$$

$$= \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\phi_i\partial\phi_j\partial\phi_k} - \mu\frac{\partial^3 H(\theta)}{\partial\theta_i\partial\theta_j\partial\theta_k}$$

From (2.9), writing $i=j$ and $j=k$ we have

$$\frac{\partial^2 H(\theta,\phi=\theta)}{\partial\theta_j\partial\phi_k} + \frac{\partial^2 H(\theta,\phi=\theta)}{\partial\phi_j\partial\phi_k} = \mu\frac{\partial^2 H(\theta)}{\partial\theta_j\partial\theta_k}.$$

Differentiating with respect to $\theta_i$

$$\frac{\partial^3 H(\theta,\phi=\theta)}{\partial\theta_i\partial\theta_j\partial\phi_k} + \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\phi_i\partial\theta_j\partial\phi_k} + \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\theta_i\partial\phi_j\partial\phi_k} + \frac{\partial^3 H(\theta,\phi=\theta)}{\partial\phi_i\partial\phi_j\partial\phi_k} = \mu\frac{\partial^3 H(\theta)}{\partial\theta_i\partial\theta_j\partial\theta_k}$$

which gives (2.14) as equivalent to (2.15).  This proves Theorem 2.2.

Let H be Shannon's entropy.  Then, an easy computation gives

$$(2.16) \quad c_{ijk} = \lambda\mu\{[\Gamma_{ijk}^{(1)}+(1-\lambda)T_{ijk}]+[\Gamma_{jki}^{(1)}+(1-\mu)T_{ijk}]+[\Gamma_{ikj}^{(1)}+(1-\mu)T_{ijk}]\}$$

where

$$(2.17) \quad \Gamma_{ijk}^{(1)} = E(\frac{\partial^2\log p_\theta}{\partial\theta_i\partial\theta_j}\frac{\partial\log p_\theta}{\partial\theta_k}), \quad T_{ijk} = E(\frac{\partial\log p_\theta}{\partial\theta_i}\frac{\partial\log p_\theta}{\partial\theta_j}\frac{\partial\log p_\theta}{\partial\theta_k}).$$

Adopting the notation of Amari for $\alpha$-connexion

$$\Gamma_{ijk}^{(\alpha)} = \Gamma_{ijk}^{(1)} + \frac{1-\alpha}{2}T_{ijk}$$

the expression (2.16) can be written

$$(2.18) \qquad c_{ijk} = \lambda\mu[\Gamma_{ijk}^{(2\lambda-1)} + \Gamma_{jki}^{(2\mu-1)} + \Gamma_{ikj}^{(2\mu-1)}].$$

When $\lambda = \mu = \frac{1}{2}$, (2.18) becomes

$$(2.19) \qquad c_{ijk} = \frac{1}{4} [\Gamma_{ijk}^{(0)} + \Gamma_{jki}^{(0)} + \Gamma_{ikj}^{(0)}].$$

Remark 1. In the definition of the Jensen difference (2.2), we used apriori probabilities $\lambda$ and $\mu$ for the two probability distributions p and q which have some relevance in population studies. But in problems of statistical inference, a symmetric version may be used by taking $\lambda = \mu = \frac{1}{2}$.

Remark 2. Throughout the discussion of this section, it was assumed that the family of probability distributions admit densities. This was done to make the computations simple. The problems could, however, be discussed in greater generality using distribution functions instead of densities.

## 3.  THE QUADRATIC ENTROPY

The quadratic entropy was introduced in Rao (1982a) as a general measure of diversity of a probability distribution over any measurable space. It is defined as a function Q: $P \rightarrow R_+$

$$(3.1) \qquad Q(p) = \int_{X \times X} K(x,y)p(x)p(y)d\nu(x)d\nu(y)$$

where $K(x,y)$ is symmetric, non-negative and conditionally negative definite, i.e.,

$$\sum_{1}^{n}\sum_{1}^{n} K(x_i,x_j)a_i a_j \leq 0$$

for any choice of $(x_1,\ldots,x_n)$ and of $(a_1,\ldots,a_n)$ such that $a_1+\ldots+a_n = 0$, with the further condition $K(x,y) = 0$ if $x = y$. As shown in Rao (1982b) and Lau (1985), the quadratic entropy is concave over $P$ and its Jensen difference has nice convexity properties which makes it an ideal measure of diversity. In view of its usefulness in statistical applications, we give explicit expressions for the quadratic differential metric and the connection coefficients associated with the quadratic entropy, in the case of the parametric family $P_\theta$.

From Theorem 2.1, the $(i,j)$-th element of the Q-information matrix is

$$(3.2) \qquad g^Q_{ij}(\theta) = -\left.\frac{\partial^2 Q(\lambda p_\theta + \mu p_\phi)}{\partial \theta_i \partial \phi_j}\right|_{\phi=\theta}.$$

Observing that

$$Q(\lambda p_\theta + \mu p_\theta) = \int K(x,y)[\lambda p(x,\theta) + \mu p(x,\phi)][\lambda p(y,\theta) + \mu p(y,\phi)]d\nu(x)d\nu(y),$$

we find the explicit expression for (3.2) as

$$(3.3) \qquad g^Q_{ij}(\theta) = -2\lambda\mu \int K(x,y) \frac{\partial p(x,\theta)}{\partial \theta_i} \frac{\partial p(y,\theta)}{\partial \theta_j} d\nu(x)\partial\nu(y)$$

$$= -2 \lambda\mu\, E[K(x,y) \frac{\partial \log p(x,\theta)}{\partial \theta_i} \frac{\partial \log p(y,\theta)}{\partial \theta_j}].$$

Using the expression (2.14), we find on carrying out the necessary computations

$$c^Q_{ijk} = -2\lambda\mu(\Gamma_{ijk} + \Gamma_{ikj} + \Gamma_{jki})$$

where

$$(3.4) \qquad \Gamma_{ijk} = \int K(x,y) \frac{\partial p(x,\theta)}{\partial \theta_k} \frac{\partial^2 p(y,\theta)}{\partial \theta_i \partial \theta_j} d\nu(x)d\nu(y).$$

It is of interest to note that the expressions (3.3) and (3.4) are invariant for transformations of both the parameters and variables.

## 4. METRICS BASED ON DIVERGENCE MEASURES

Burbea and Rao (1982a, 1982b), Burbea (1984) and Eguchi (1984) have considered metrics arising out of a variety of divergence measures between probability distributions. A typical divergence measure is of the form

$$(4.1) \qquad D_F(p_\theta, p_\phi) = \int_X F[p(x,\theta), p(x,\phi)] d\nu(x)$$

where F satisfies the following conditions:

(i) $F(\cdot, \cdot)$ is a $C^3$-function on $R_+ \times R_+$,

(ii) $F(x, \cdot)$ is strictly convex on $R_+$ for every $x \in R_+$,

(iii) $F(x,x) = 0$ for every $x \in R_+$,

(iv) $\dfrac{\partial F(x, y = x)}{\partial y} = $ constant for every $x \in R_+$.

Let us consider the expansion

$$(4.2) \qquad D_F(p_\theta, p_{\theta + d\theta}) = \frac{1}{2!} \Sigma\Sigma g^F_{ij}(\theta) d\theta_i d\theta_j + \frac{1}{3!} c^F_{ijk}(\theta) d\theta_i d\theta_j d\theta_k + \ldots$$

and obtain explicit expressions for $g^F_{ij}$ and $c^F_{ijk}$.

Theorem 4.1. Let

$$F_1(x,y) = \frac{\partial F(x,y)}{\partial x} , \quad F_2(x,y) = \frac{\partial F(x,y)}{\partial y}$$

$$F_{11} = \frac{\partial^2 F(x,y)}{\partial x^2} , \quad F_{12} = \frac{\partial^2 F(x,y)}{\partial x \partial y} , \quad F_{22} = \frac{\partial^2 F(x,y)}{\partial y^2}$$

$$F_{222} = \frac{\partial^3 F(x,y)}{\partial y^3} .$$

Then

$$(i) \quad g^F_{ij}(\theta) = \int F_{22}[p_\theta, p_\theta] \frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x)$$

$$= -\int F_{12}[p_\theta, p_\theta] \frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x).$$

(ii) $c_{ijk}^F = \int F_{222}[p_\theta, p_\theta] \frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} \frac{\partial p_\theta}{\partial \theta_k} d\nu(x)$

$$+ \int F_{22}[p_\theta, p_\theta][\frac{\partial^2 p_\theta}{\partial \theta_i \partial \theta_j} \frac{\partial p_\theta}{\partial \theta_k} + \frac{\partial^2 p_\theta}{\partial \theta_i \partial \theta_k} \frac{\partial p_\theta}{\partial \theta_j} + \frac{\partial^2 p_\theta}{\partial \theta_j \partial \theta_k} \frac{\partial p_\theta}{\partial \theta_i}]d\nu(x).$$

The results are established by straight forward computations.

Let us consider the directed divergence measure of Csiszár (1967), which plays an important role in problems of statistical inference,

(4.3)  $$D(p_\theta, p_\phi) = \int p(x,\theta) \, f(\frac{p(x,\phi)}{p(x,\theta)}) d\nu(x)$$

where f is a convex function.  In this case

(4.4)  $$g_{ij}^f(\theta) = \frac{\partial^2 D}{\partial \phi_i \partial \phi_j} \Big|_{\phi=\theta}$$

$$= f''(1) \int \frac{1}{p} \frac{\partial p}{\partial \theta_i} \frac{\partial p}{\partial \theta_j} d\nu(x) = f''(1) g_{ij}(\theta)$$

where $g_{ij}$ are the elements of Fisher's information matirx.  Thus a wide class of invariant divergence measures provide the same informative geometry on the parameter manifold.  Further,

$$c_{ijk}^f(\theta) = \frac{\partial^3 D}{\partial \phi_i \partial \phi_j \partial \phi_k} \Big|_{\phi=\theta}$$

$$= f''(1)[\Gamma_{ijk}^{(1)} + \Gamma_{ikj}^{(1)} + \Gamma_{jki}^{(1)}] + (f'''(1) + 3f''(1))T_{ijk}$$

where $\Gamma_{ijk}^{(1)}$ and $T_{ijk}$ are as defined in (2.17).

If f is a convex function, then

$$f^*(u) = uf(\frac{1}{u})$$

is also convex, and the measure (4.3) associated with $f + f^*$ is

$$(4.5) \qquad D^*(p_\theta, p_\phi) = \int [p_\theta f(\frac{p_\phi}{p_\theta}) + p_\phi f(\frac{p_\theta}{p_\phi})] d\nu(x)$$

which is symmetric in $\theta$ and $\phi$. However, we may define (4.5) as a symmetric divergence measure without requiring f to be a convex function but satisfying the condition that $xf(x^{-1}) + f(x)$ is non-negative on $R_+$. In such a case

$$g_{ij}^f(\theta) = 2f''(1)g_{ij}(\theta)$$

$$c_{ijk}^f(\theta) = f''(1)[\Gamma_{ijk}^{(1)} + \Gamma_{ikj}^{(1)} + \Gamma_{jki}^{(1)}] + f'''(1)T_{ijk}$$

## 5. OTHER DIVERGENCE MEASURES

In the last section, we considered the f-divergence measure which led to the Fisher information metric. A special case of this measure is the city block distance, or the overlap distance (see Rao, 1948, 1982a),

$$(5.1) \qquad D_0(p_\theta, p_\phi) = \int |p(x,\theta) - p(x,\phi)| d\nu(x)$$

obtained by choosing $f(x) = 1 - \min(x,1)$, which admits a direct interpretation in terms of errors of classification in discrimination problems. However, this is not a smooth function and no formula of the type (4.7) is avialable to determine the coefficients of the differential metric. But in some cases, it may turn out that

$$D_0(p_\theta, p_\phi) = D_0(\theta, \phi)$$

is a smooth function of $\theta$ and $\phi$ in which case

$$(5.2) \qquad g_{ij} = \frac{\partial^2 D_0(\theta, \phi=\theta)}{\partial \phi_i \partial \phi_j}.$$

In the case when $p(x,\theta)$ is a p-variate normal density with mean $\mu$ and fixed variance covariance matrix $\Sigma$, the coefficient (5.2) can be easily computed to be proportional to $\sigma^{ij}$, the (i,j)-th element of $\Sigma^{-1}$, which is indeed the (i,j)-th

element of the Fisher information matrix. It would be of interest to investigate the nature of the metric induced by (5.1) in the general case.

Let $p(x,\theta)$ be the density of a uniform distribution in the interval $[0,\theta]$. Then it is seen that

$$(5.3) \qquad D_0(\theta,\phi) = 2(1 - \frac{\theta}{\phi}) \quad \text{if } \theta \leq \phi$$
$$= 2(1 - \frac{\phi}{\theta}) \quad \text{if } \theta > \phi.$$

Although this is not a differentiable function, it is seen that

$$ds^2 = 4\frac{d\theta^2}{\theta^2}$$

is the metric associated with (5.3).

Another general divergence measure which has some practical applications is

$$D_\psi(p_\theta, p_\phi) = \int [\psi(p_\theta) - \psi(p_\phi)]^2 d\nu(x)$$

which is indeed a smooth function if $\psi$ is so. In this case

$$g_{ij}^\psi(\theta) = 2\int [\psi'(p_\theta)]^2 \frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x)$$

$$c_{ijk}^\psi(\theta) = 6 \int \psi'(p_\theta)\psi''(p_\theta)\frac{\partial p_\theta}{\partial \theta_i} \frac{\partial p_\theta}{\partial \theta_j} \frac{\partial p_\theta}{\partial \theta_k} d\nu(x)$$
$$+ 2 \int [\psi'(p_\theta)]^2 (\frac{\partial^2 p_\theta}{\partial \theta_i \partial \theta_j} \frac{\partial p_\theta}{\partial \theta_k} + \frac{\partial^2 p_\theta}{\partial \theta_i \partial \theta_k} \frac{\partial p_\theta}{\partial \theta_j} + \frac{\partial^2 p_\theta}{\partial \theta_j \partial \theta_k}) \frac{\partial p_\theta}{\partial \theta_i} d\nu(x)$$

Another measure of interest is the cross entropy introduced in Rao and Nayak (1985). If H is any entropy function, then the cross entropy of $p_\phi$ with respect to $p_\theta$ was defined as

$$(5.4) \quad D(p_\theta, p_\phi) = H(p_\phi) - H(p_\theta) - \lim_{\lambda \to 0} \frac{H[p_\phi + \lambda(p_\theta - p_\phi)] - H(p_\phi)}{\lambda}.$$

Let

$$H(p) = -\int h(p) d\nu(x)$$

as chosen in (2.10). Then (5.4) reduces to

$$D(p_\theta, p_\phi) = -\int h(p_\phi) d\nu(x) - \int h'(p_\phi)(p_\theta - p_\phi) d\nu(x) + \int h(p_\theta) d\nu(x).$$

Then

$$g_{ij}^h = \int h''(p_\theta) \frac{\partial p_\theta}{\partial \theta_j} \frac{\partial p_\theta}{\partial \theta_j} d\nu(x)$$

which is the same as the h-entropy information matrix derived in (2.10), apart from a constant. Similarly

$$c_{ijk}^h = \Gamma_{ijk}^{(1)} + \Gamma_{ikj}^{(1)} + \Gamma_{jki}^{(1)} + T_{ijk}$$

where

$$\Gamma_{ijk}^{(1)} = E\{p_\theta h''(p_\theta) \frac{\partial^2 \log p_\theta}{\partial \theta_i \partial \theta_j} \frac{\partial^2 \log p_\theta}{\partial \theta_k}\}$$

$$T_{ijk} = E\{[3p_\theta h''(p_\theta) + 2p_\theta^2 h'''(p_\theta)] \frac{\partial \log p_\theta}{\partial \theta_i} \frac{\partial \log p_\theta}{\partial \theta_j} \frac{\partial \log p_\theta}{\partial \theta_k}\}.$$

## 6. GEODESIC DISTANCES

In Rao (1945) it was suggested that the information metric could be used to obtain the geodesic distances between probability distributions. Given any quadratic differential metric

$$(6.1) \quad ds^2 = \Sigma\Sigma \, g_{ij}(\theta) d\theta_i d\theta_j$$

where the matrix $(g_{ij})$ is positive definite, the geodesic curve $\theta = \theta(t)$ can be determined from the Euler-Lagrange equations

$$(6.2.) \qquad \sum_{1}^{n} g_{ik} \ddot{\theta}_i + \sum_{1}^{n}\sum_{1}^{n} \Gamma_{ijk} \dot{\theta}_i \dot{\theta}_j = 0, \quad k = 1, \ldots, n$$

and from the boundary conditions

$$\theta(t_1) = \theta, \quad \theta(t_2) = \phi.$$

In (6.2), the quantity

$$(6.3) \qquad \Gamma_{ijk} = \frac{1}{2}[\frac{\partial}{\partial\theta_i} g_{jk} + \frac{\partial}{\partial\theta_j} g_{ki} - \frac{\partial}{\partial\theta_k} g_{ij}]$$

and is known as the "Christoffel symbol of the first kind".

By definition of the geodesic curve $\theta = \theta(t)$, its tangent vector $\dot{\theta} = \dot{\theta}(t)$ is of constant length with respect to the metric $ds^2$. Thus

$$(6.4) \qquad \sum_{1}^{n}\sum_{1}^{n} g_{ij} \dot{\theta}_i \dot{\theta}_j = \text{constant}.$$

The constant may be chosen to be of value 1 when the curve parameter $t$ is the arc length parameter $s$, $0 \le s \le s_o$, with $\theta(0) = \theta$, $\theta(s_0) = \phi$ and $s_0 = g(\theta,\phi)$ is the geodesic distance between $\theta$ and $\phi$.

Aitkinson and Mitchell (1981) describe two other methods of deriving geodesic distances starting from a given differential metric. The distances obtained by these authors in various cases are given below. In each case we give the probability function $p(x,\theta)$ and the associated geodesic distance of $(\theta,\phi)$ based on the Fisher information metric.

(1) <u>Poisson distribution</u>

$$p(x,\theta) = e^{-\theta} \theta^x/x!, \quad x = 0,1,\ldots$$

$$g(\theta,\phi) = 2|\sqrt{\theta} - \sqrt{\phi}|$$

(2) **Binomial distribution (n fixed)**

$$p(x,\theta) = \binom{n}{x}\theta^x(1-\theta)^{n-x}, \quad x = 0,1,\ldots,n$$

$$g(\theta,\phi) = 2\sqrt{n}\,|\sin^{-1}\sqrt{\theta} - \sin^{-1}\sqrt{\phi}|$$

$$= 2\sqrt{n}\,\cos^{-1}[\sqrt{\theta\phi} + \sqrt{(1-\theta)(1-\phi)}\,].$$

(3) **Exponential distribution**

$$p(x,\theta) = \theta e^{-x\theta}, \quad x \geq 0$$

$$g(\theta,\phi) = |\log\theta - \log\phi|.$$

(4) **Gamma distribution (n fixed)**

$$p(x,\theta) = \theta^n[\Gamma(n)]^{-1}x^{n-1}e^{-x\theta}, \quad x \geq 0$$

$$g(\theta,\phi) = \sqrt{n}\,|\log\theta - \log\phi|$$

(5) **Normal distribution**

$$p(x,\mu,\sigma_0^2) = N(\mu,\sigma_0^2;x), \quad \sigma_0 \text{ fixed}$$

$$g(\mu_1,\mu_2) = |\mu_1 - \mu_2|/\sigma_0$$

(6) **Normal distribution**

$$p(x,\mu_0,\sigma^2) = N(\mu_0,\sigma^2;x), \quad \mu_0 \text{ fixed}$$

$$g(\sigma_1^2,\sigma_2^2) = \sqrt{2}\,|\log\sigma_1 - \log\sigma_2|$$

(7) **Normal distribution**

$$p(x,\mu;\sigma^2) = N(\mu,\sigma^2;x), \quad \mu \text{ and } \sigma \text{ both variable.}$$

The information metric in this case is

(6.5)
$$ds^2 = \frac{d\mu^2}{d\sigma^2} + \frac{2d\sigma^2}{\sigma^2}$$

and the geodesic distance is

(6.6)
$$g(\mu_1,\sigma_1;\mu_2,\sigma_2) = \sqrt{2}\,\left|\log\frac{1+\delta(1,2)}{1-\delta(1,2)}\right|$$

$$= 2\sqrt{2}\,\tanh^{-1}\delta(1,2)$$

where $\delta_{12}$ is the positive square root of

$$\frac{(\mu_1-\mu_2)^2 + 2(\sigma_1-\sigma_2)^2}{(\mu_1-\mu_2)^2 + 2(\sigma_1+\sigma_2)^2} \ .$$

The explicit form (6.6) is given in Burbea and Rao (1982a). From (6.6)

$$g(\mu,\sigma_1^2;\mu,\sigma_2^2) = \sqrt{2} \, |\log \sigma_1 - \log \sigma_2|$$

which agrees with result (6). However, $g(\mu_1,\sigma^2;\mu_2,\sigma^2)$ does not reduce to result (7) since $\sigma$ = constant is not a geodesic curve with respect to the metric (6.5)

(8) <u>Multivariate normal distribution</u>

$N_p(\mu,\Sigma;x)$, $\Sigma$ fixed

$$g(\mu_1,\mu_2) = (\mu_1-\mu_2)'\Sigma^{-1}(\mu_1-\mu_2)$$

which is Mahalanobis distance.

(9) <u>Multivariate normal distribution</u>

$N(\mu,\Sigma;x)$, $\mu$ fixed

$$g(\Sigma_1,\Sigma_2) = 2^{-1} \sum_{1}^{p} (\log \lambda_i)^2$$

where $0 < \lambda_1 \le \ldots \le \lambda_p$ are the roots of the determinantal equation $|\Sigma_2 - \lambda\Sigma_1| = 0$. The above explicit form is due to S.T. Jensen as mentioned in Atkinson and Mitchell (1981).

(10) <u>Negative binomial distribution</u>

$p(x,\theta) = [x!\Gamma(r)]^{-1}\Gamma(x+r)\theta^x(1-\theta)^r$, $r$ fixed

$$g(\theta,\phi) = 2\sqrt{r} \, \cosh^{-1} \frac{1 - \sqrt{\theta\phi}}{\sqrt{(1-\theta)(1-\phi)}}$$

$$= 2\sqrt{r} \, \log \frac{1-\sqrt{\theta\phi}+|\sqrt{\theta} - \sqrt{\phi}|}{\sqrt{(1-\theta)(1-\phi)}}$$

This computation is due to Oller and Cuadras (1985).

(11)  **Multinomial distribution**

$$p(n_1,\ldots,n_k; \pi_1,\ldots,\pi_k) = \frac{n!}{n_1!\ldots n_k!}\, \pi_1^{n_1}\ldots\pi_k^{n_k}\ , \quad n \text{ fixed.}$$

Let $\underset{\sim}{\pi}_1 = (\pi_{11},\ldots,\pi_{k1})$ and $\underset{\sim}{\pi}_2 = (\pi_{12},\ldots,\pi_{k2})$.  Then

$$g(\underset{\sim}{\pi}_1,\underset{\sim}{\pi}_2) = 2\sqrt{n}\ \cos^{-1}(\sum_1^k \sqrt{\pi_{i1}\pi_{i2}}\ )$$

The above computation was originally done by Rao (1945), but an easier method

of derivation is given by Atkinson and Mitchell (1981).

Recently Burbea (1984) obtained geodesic distances in the case of independent Poisson and Normal distributions which are given below.

(12)  **Independent Poisson distributions**

$$p(x_1,\ldots,x_n;\theta_1,\ldots,\theta_n) = \prod_1^n e^{-\theta_i}\ \frac{\theta_i^{x_i}}{x_i!}$$

$$g(\theta_1,\ldots,\theta_n;\phi_1,\ldots,\phi_n) = 2[\sum_1^n(\sqrt{\theta_i} - \sqrt{\phi_i}\ )^2]^{1/2}$$

(13)  **Independent Normal distributions**

$$N(x;\mu_1,\sigma_1^2)\ldots N(x_n;\mu_n,\sigma_n^2)$$

$$g[(\mu_{11}\sigma_{11}^2),\ldots,(\mu_{n1},\sigma_{n1}^2);(\mu_{12},\sigma_{12}^2),\ldots,(\mu_{n2},\sigma_{n2}^2)]$$

$$= \sqrt{2}\ [\sum_{k=1}^n \log^2 \frac{1+\delta_k(1,2)}{1-\delta_k(1,2)}]^{1/2}$$

where $\delta_k(1,2)$ is the positive square root of

$$\frac{(\mu_{k1}-\mu_{k2})^2 + 2(\sigma_{k1}-\sigma_{k2})^2}{(\mu_{k1}-\mu_{k2})^2 + 2(\sigma_{k1}+\sigma_{k2})^2}\ .$$

(14)  **Multivariate elliptic distributions**

$$p(x|\underset{\sim}{\mu},\Sigma) = |\Sigma|^{-1/2}h[(x-\mu)'\Sigma^{-1}(x-\mu)],$$

for some function h, and $\Sigma$ is fixed

$$g(\underset{\sim}{\mu}_1, \underset{\sim}{\mu}_2) = c_h (\underset{\sim}{\mu}_1 - \underset{\sim}{\mu}_2)' \Sigma^{-1} (\underset{\sim}{\mu}_1 - \underset{\sim}{\mu}_2)$$

where $c_h$ is a constant, which is essentially Mahalanobis distance. This result is due to Mitchell and Krzanowski (1985).

The use of the $c_{ijk}$ coefficients defined in (2.4) and (4.2) in the discussion of statistical problems will be considered in a future communication.

## 7. BIBLIOGRAPHY

[1] Amari, S.I. (1982). Differential geometry of curved exponential families-curvature and information loss. Ann. Statist. 10, 357-385.

[2] Amari, S.I. (1983). A foundation of information geometry. Electronics and Communications in Japan 66-A, 1-10.

[3] Atkinson, C. and MItchell, A.F.S. (1981). Rao;s distance measure. Sankhya 43, 345-365.

[4] Burbea, J. (1984). Informative geometry in probability spaces. Tech. Report No. 84-52, Center for Multivariate Analysis, University of Pittsburgh.

[5] Burbea, J. and Rao, C. Radhakrishna (1982a). Entropy differential metric, distance and divergence measures in probability spaces: a unified approach. J. Multivariate Anal. 12, 576-596.

[6] Burbea, J. and Rao, C. Radhakrishna (1982b). Differential metrics in probability. 3, 115-132.

[7] Čencov, N.N. (1982). Statistical decision rules and optimal inference. Transactions of Mathematical Monographs 53, Amer. Math. Soc., Providence.

[8] Csiszár, I. (1967). Information-type measures of difference of probability distributions and indirect observations. Studia Scientiarum Mathematicarum Hungrica 2, 299-318.

[9] Efron, B. (1975). Defining the curvature of a statistical problem (with applications to second order efficiency, with discussion). Ann. Statist. 3, 1189-1217.

[10] Efron, B. (1982). Maximum likelihood decision theory. Ann. Statist. 10, 340-356.

[11] Eguchi, Shinto (1983). Second order efficienty of minimum contrast estimators in a curved exponential family. Ann. Statist. 11, 793-803.

[12] Eguchi, Shinto (1984). A differential geometric approach to statistical inference on the basis of contrast functionals. Tech. Report No. 136, Hiroshima University, Hiroshima, Japan.

[13] Lau, Ka-Sing (1985). Characterization of Rao's quadratic entropy.
     Sankhya A (to appear).

[14] Mitchell, A.F.S. and Krzanowski, W.J. (1985). The Mahalanobis distance
     and elliptic distributions. (To appear in Biometrika).

[15] Nei, M. (1978). The theory of genetic distance and evolution of human
     races. Japan J. Human Genet. 23, 341-369.

[16] Oller, J.M. and Cuadras, C.M. (1985). Rao's distance for negative multi-
     nomial distributions. Sankhya 47, 75-83.

[17] Rao, C. Radhakrishna (1945). Information and accuracy attainable in the
     estimation of statistical parameters. Bull. Calcutta Math. Soc. 37,
     81-91.

[18] Rao, C. Radhakrishna (1948). The utilization of multiple measurements in
     problems of biological classification (with discussion). J. Roy. Statist.
     Soc. B10, 159-203.

[19] Rao, C. Radhakrishna (1949). On the distance between two populations.
     Sankhya 9, 246-248.

[20] Rao, C. Radhakrishna (1954). On the use and interpretation of distance
     functions in statistics. Bull. Inst. Inter. Statist. 34, 90-100.

[21] Rao, C. Radhakrishna (1962). Efficient estimates and optimum infernece
     procedures in large samples (with discussion). J. Roy. Statist. Soc. B,
     24, 46-72.

[22] Rao, C. Radhakrishna (1973). Linear Statistical Inference and its Applications.
     (Second edition) Wiley, New York.

[23] Rao, C. Radhakrishna (1982a). Diversity and dissimilarity coefficients:
     a unified approach. J. Theoret. Pop. Biology 21, 24-43.

[24] Rao, C. Radhakrishna (1982b). Diversity: its measurement, decomposition,
     apportionment and analysis. Sankhya A, 44, 1-22.

[25] Rao, C. Radhakrishna and Nayak, T.K. (1985). Cross entropy, dissimilarity
     measures and characterizations of quadratic entropy. IEEE Trans Information
     Theory (to appear).

| REPORT DOCUMENTATION PAGE | READ INSTRUCTIONS<br>BEFORE COMPLETING FORM |
|---|---|

| 1. REPORT NUMBER<br>AFOSR-TR- 4 | 2. GOVT ACCESSION NO.<br>AD-A160 301 | 3. RECIPIENT'S CATALOG NUMBER |
|---|---|---|

| 4. TITLE (and Subtitle)<br>Differential Metrics in Probability Spaces Based on Entropy and Divergence Measures | 5. TYPE OF REPORT & PERIOD COVERED<br>Technical - April, 1985 |
|---|---|
| | 6. PERFORMING ORG. REPORT NUMBER<br>85-08 |

| 7. AUTHOR(s)<br>C. Radhakrishna Rao | 8. CONTRACT OR GRANT NUMBER(s)<br>F-49620-85-C-0008 |
|---|---|

| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>Center for Multivariate Analysis<br>Room 515, Thackeray Hall, University of Pittsburgh<br>Pittsburgh, Pennsylvania 15260 | 10. PROGRAM ELEMENT, PROJECT, TASK<br>AREA & WORK UNIT NUMBERS<br>61102F<br>2304/A5 |
|---|---|

| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>Air Force Office of Scientific Research<br>Department of the Air Force<br>Bolling Air Force Base, DC 20332 | 12. REPORT DATE<br>April, 1985 |
|---|---|
| | 13. NUMBER OF PAGES<br>27 |

| 14. MONITORING AGENCY NAME & ADDRESS(If different from Controlling Office) | 15. SECURITY CLASS. (of this report)<br>Unclassified |
|---|---|
| | 15a. DECLASSIFICATION/DOWNGRADING<br>SCHEDULE |

**16. DISTRIBUTION STATEMENT (of this Report)**

Approved for public release; distribution unlimited

**17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)**

**18. SUPPLEMENTARY NOTES**

**19. KEY WORDS (Continue on reverse side if necessary and identify by block number)**

cross entropy, f-divergence, geodesic distance, information matrix, Jensen difference, quadratic entropy

**20. ABSTRACT (Continue on reverse side if necessary and identify by block number)**
In this paper are discussed some general methods of metrizing probability spaces through the introduction of a quadratic differential metric in the parameter manifold of a set of probability distributions. These methods extend the investigation made in Rao (1945) where the Fisher information matrix was used to construct the metric, and the geodesic distance was suggested as a measure of dissimilarity between probability distributions.
(Continued)

DD FORM 1473
1 JAN 73

Block #20, "Abstract", Continued

The basic approach in the present paper is first to construct a divergence
or a dissimilarity measure between any two probability distributions, and
use it to derive a differential metric by considering two distributions
whose characterizing parameters are close to each other.  One measure of
divergence considered is the Jensen difference based on an entropy functional
as defined in Rao (1982a).  Another is the f-divergence measure studied by
Csiszár (1967).  The latter class leads to the differential metric based on
the Fisher information matrix.  The geodesic distances based on this metric
computed by various authors are listed.

# END

# FILMED

11-85

# DTIC